

10-02-00

A

09/29/00
10715 U.S. PTO

UTILITY PATENT APPLICATION TRANSMITTAL		Attorney Docket No.	40921/205584
		First Named Inventor	Keith Glidewell
		Title	APPARATUS AND METHOD FOR PROCESS DISPATCHING BETWEEN INDIVIDUAL PROCESSORS OF A MULTI-PROCESSOR SYSTEM
		Express Mail Label No.	EJ841202327US

APPLICATION ELEMENTS
See MPEP chapter 600 concerning utility patent application contents

ADDRESS TO: Assistant Commissioner for Patents
Box Patent Application
Washington, D.C. 20231

1. ☒ Specification Total Pages 16
(preferred arrangement as set forth below)
 - Descriptive title of the Invention
 - Cross References to Related Applications
 - Statement Regarding Fed sponsored R & D
 - Reference to Microfiche Appendix
 - Background of the Invention
 - Brief Summary of the Invention
 - Brief Description of the Drawings (if filed)
 - Detailed Description
 - Claim(s)
 - Abstract of the Disclosure
2. ☒ Drawing(s) (35 USC 113) Total Sheets 7
3. Oath or Declaration Total Pages 2
 - a. ☒ Newly executed (original or copy)
 - b. ☐ Copy from a prior application (37 CFR 1.63(d)) (for continuation/divisional with box 17 completed) [Note Box 4 below]
 - i. ☐ DELETION OF INVENTOR(S)
Signed statement attached deleting inventor(s) named in the prior application, See 37 CFR 1.63(d)(2) and 1.33(b)
4. ☐ Incorporation By Reference (usable if Box 3b is checked)
The entire disclosure of the prior application, which a copy of the oath or declaration is supplied under Box 3b, is considered as being part of the disclosure of the accompanying application and is hereby incorporated by reference therein

5. ☐ Microfiche Computer Program (Appendix)
6. Nucleotide and/or Amino Acid Sequence Submission (if applicable, all necessary)
 - a. ☐ Computer Readable Copy
 - b. ☐ Paper Copy (identical to computer copy)
 - c. ☐ Statement verifying identify of above copies

ACCOMPANYING APPLICATION PARTS


7. ☒ Assignment Papers (cover sheet & document(s))
 8. ☐ 37 CFR 3.73(b) Statement ☒ Power of Attorney (when there is an assignee)
 9. ☐ English Translation Document (if applicable)
 10. ☐ Information Disclosure Statement (IDS)/PTO-1449 ☐ Copies of IDS Citations
 11. ☐ Preliminary Amendment
 12. ☒ Return Receipt Postcard
 13. ☐ Small Entity Statement ☐ Statement filed in prior application, Status still proper and desired
 14. ☐ Certified Copy of Priority Document(s) (If foreign priority is claimed)
 15. ☒ Other:
 1. Certificate of Express Mail No. EJ841202327US
 2. Associate Power of Attorney
- Check in the amount of \$730.00 The Commissioner is hereby authorized to credit overpayments or charge any additional fees under 37 CFR 1.116 and 1.117 to Deposit Account No. 16-1435

16. If a **CONTINUING APPLICATION**, check appropriate box and supply the requisite information below and in a preliminary amendment.
☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application
 Prior application information: Examiner: Group/Art Unit:

17. FEE CALCULATIONS

CLAIMS	For	Number Filed	Extra	Rate	Calculations
	Total Claims	20 - 20 =	0 x	\$18 =	\$ 0.00
	Indep. Claims	2 - 3 =	0 x	\$78 =	\$ 0.00
	Multiple Dependent Claims (if applicable)		+	\$260 =	\$ 0.00
				Basic Fee (37 CFR 1.16)	\$ 690.00
				Total Calculations	\$ 690.00
				Reduced by 50% for filing small entity (Note 37 CFR 1.9, 1.27, 1.28).	\$
	Fee for recording the enclosed assignment (37 CFR 1.21(h)). The assignment must be accompanied by an appropriate cover sheet (37 CFR 3.28, 3.31) \$40 per property				\$ 40.00
	TOTAL FEES SUBMITTED				\$ 730.00

18. CORRESPONDENCE ADDRESS

Name	A. Jose Cortina, Esquire KILPATRICK STOCKTON LLP		
Address	3737 Glenwood Avenue Suite 400		
City	Raleigh	State	NC
Country	U.S.A.	Telephone	919-420-1700
		Fax	919-420-1800
Name (Print/Type)	A. Jose Cortina	Registration No (Attorney/Agent)	29,733
Signature		Date	9/29/2000

10715 U.S. PTO
09/29/00

09/29/00

EXPRESS MAIL CERTIFICATE

"Express Mail" mailing label number: EJ841202327US

Applicant: Keith Glidewell

Title: Apparatus and Method for Process Dispatching
Between Individual Processors of a
Multi-Processor System

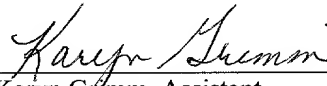
Date of Deposit Herewith

Type of Document(s) Utility Patent Application Transmittal (1 page);
Specification and Abstract; 16 pages
Including 20 claims;
Drawings (7 sheets);
Declaration and Power of Attorney (2 pages);
Associate Power of Attorney (1 page);
Assignment of Invention and Cover Sheet (4 pages);
Check #19189 in the amount of \$730.00;
Express Mail Certificate;
Return postcard.

Serial No.: Unassigned

Filing Date: Herewith

I hereby certify that the documents identified above are being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on September 29, 2000 and is addressed to the Assistant Commissioner for Patents, Box Patent Application, Washington, D.C. 20231.


Karyn Grimm, Assistant

APPARATUS AND METHOD FOR PROCESS DISPATCHING
BETWEEN INDIVIDUAL PROCESSORS OF A
MULTI-PROCESSOR SYSTEM

5

BACKGROUND OF THE INVENTION

1. Field of the Invention

The present invention relates generally to computer operating systems and more particularly to a method for dispatching processes between
10 individual processors of a multi-processor system.

2. Description of the Prior Art

In a multi-processing system on which symmetric multi-processing is conducted, each processor has equal access to the memory and input/output resources. The problem with such systems is that as you add processors to the
15 BUS, the BUS becomes saturated and it becomes impossible to add more processors.

A prior art solution to that problem was to extend symmetric multi-processing with a technique known as Cache Coherent Non-Uniform Memory Access, i.e., ccNUMA. In such systems, the processors do not have to do anything
20 extraordinary to access memory. Normal reads and writes are issued, and caching is all handled in hardware. By the term "Coherent" is meant that if processor "A" reads a memory location and processor "B" reads a memory location, both processors see the same data at that memory location. If one processor writes to that location, the other will see the data written. The fact that access is described
25 as being non-uniform refers to the fact that memory is physically and electrically closer to one processor than to another processor. So if memory is accessed from one processor, and it is local, the access is much faster than if the memory is accessed from a remote processor.

Modern multi-processing systems can have eight or more individual
30 processors sharing processing tasks. Many such systems incorporate caches that

are shared by a subset of the system's processors. Typically the systems are implemented in blocks, where each block is a symmetric multi-processing system consisting of four processors, and current systems have implemented up to sixteen blocks for a total arrangement of up to sixty-four processors.

5 In such a system, each block with a reference to data to be accessed is put on a sharing chain which is a link list, maintained by the ccNUMA hardware. Each additional block that has a reference to the data will have a link on the sharing chain. When a write operation occurs, the sharing chain has to be torn down, i.e., invalidated, because only one block can be writing to a memory
10 location at a time. If a memory location is being read, multiple blocks can read the same memory location with low latency. As additional blocks read the data, the sharing chain blocks can be built up with blocks that are reading the data and a copy of the data is cached local to each of these blocks.

 In such a system, best access time occurs if a block is accessing local
15 memory, i.e., memory that is from the same symmetric multi-processing block. The next best access time occurs if it is memory that is from another block, but is in the processor block's far memory cache. The next best scenario, is if a processor writes to a far memory location that was just written on another block. The reason that this does not cause excessive latency is because there is only one
20 block on the sharing chain, since the write tore down the sharing chain. Thus, if the processor has to write to the memory location, it just has to invalidate one block. The worst case scenario is where a long sharing chain has to be torn down. If eight blocks have all read the memory, and a processor wishes to write to the memory, an indication has to be sent to every one of the eight blocks indicating
25 that the copy of the memory location is invalid and can no longer be used.

 One prior art system involved adding additional processors to a NUMA system. When performance was tested, it was unexpectedly uncovered that as processors were added, through-put actually declined and the number of transactions that could be processed in a minute declined. A confusing part about
30 the decline was that an analysis of the system revealed that it was mostly idle, but that there was a significant amount of scalable coherent interconnect traffic. The

“scalable coherent interconnect” is typically the NUMA backplane. Further analysis revealed that the problem resulted from processors which were not running a user process, i.e., processors which were technically idle, but as a result of the idle state were actually spinning constantly searching for tasks to process.

5 An example of the way tasks are arranged in such a system are described in greater detail in U.S. Patent 5,745,778 which describes a method of operation of a multi-processor data processing system using a method of organizing and scheduling threads. Specifically, in that system, “threads,” which are programming constructs that facilitate efficient control of numerous asynchronous
10 tasks are assigned priorities in a limited portion of a range. The disclosure of the noted patent describes how thread groups and threads can be assigned priority and scheduled globally to compete for central processing unit, i.e., processor or CPU, resources available anywhere in the system.

 In the above-described system in which additional processors were
15 added, it was discovered that data structures used by the idle processor were not allocated from local memory, and as the processors were idle in their own data structures, they would actually be reading and writing memory from a far locale. The other problem uncovered was that the processors were searching for work too aggressively, and upon determining that there was no work to do on their own
20 queue, would start to examine other lists for other processors in other locales. As such, the processor would immediately poach a process from another processor, even if that other processor or system was suddenly going to become idle. This would cause moving all of the process’ data and cache footprint to the idle processor, resulting in reduced through-put.

25 For purposes of this disclosure, it should be noted that by the term “poaching” is meant taking a job from another processor or locale’s ready list. By “locale” is meant a symmetric multi-processor system which is implemented within a NUMA system as a block and is made up of four individual processors, each with their own cache. Similarly central processing unit, CPU and processor
30 are used interchangeably herein, and are used to refer to individual processors

arranged in symmetric multi-processing blocks, i.e., SMP deployed in a NUMA system.

Accordingly, in accordance with the invention, the delays associated with such poaching in a large multi-processor system are avoided by the system and method described herein.

SUMMARY OF THE INVENTION

There is disclosed a method of operation of a multi-processor data processing system which attempts to keep all of the processors on the multi-processor Non-Uniform Memory Architecture system performing productive work.

The method involves allocating resources in a plurality of processors system, with each processor having a cache associated therewith. When at least one processor of the plurality of processors is idle, it is determined if at least one other processor of the plurality of processors is not idle. If the processor which is not idle remains not idle for a predetermined period of time, the idle processor poaches a process on the queue of the non-idle processor to be run by the processor which is idle.

In a more specific aspect, if more than one processor is idle, poaching occurs with the idle processor which is electrically closest to the non-idle processor. The further away in terms of electrical connection proximity that a non-idle processor is to an idle processor, the longer the time period that the non-idle processor is allowed to remain non-idle.

In an alternative aspect, the invention relates to a data processing system for allocating resources for simultaneously executing a plurality of processing tasks. The system includes a plurality of processors, each having a cache associated therewith. A timer is associated with each processor for timing from the beginning of running a process from the processor's queue, the duration of time the process is run. During this time the processor is considered as being non-idle, a detector determines the duration of time any one processor is not idle, and

the processors are configured such that they can poach a process from a non-idle processor when the duration of time during which the non-idle processor is running exceeds a predetermined amount.

5 In a more specific aspect, the system is configured for allowing an idle processor connected electrically closest to a non-idle processor to have priority in poaching a process in the event there are more than one non-idle processors. The system is further configured to allow the time period during which a processor is allowed to remain non-idle, determined in accordance with the proximity in electrical connection between the non-idle processor and an idle processor, in
10 which the greater the connection distance, the greater the predetermined amount of time allowed.

Other features of the invention will be understood by those of ordinary skill in the art after referring to the detailed description of the preferred embodiment and drawings.

15 BRIEF DESCRIPTION OF THE INVENTION

Fig. 1 is an overview of a multi-processor data processing system.

Fig. 2A-2E shows the flow of how a process is stolen or poached by an idle processor when the processor on which it is scheduled on the queue is too busy to run the process.

20 Fig. 3 shows the flow of how the job processes and relatives table is allocated such that the shortest relative time out for a job process is arranged at the top of the table, and showing how time periods are set in accordance with differences in distances between processors.

DESCRIPTION OF THE PREFERRED EMBODIMENT

25 1. System Overview

Referring to Fig. 1, an overview of a multi-processing data processing system 150 is depicted. For clarity and ease of presentation, an eight-processor

system has been depicted. As will be readily appreciated by those of ordinary skill in the art, the invention is applicable to multi-processor systems having other numbers of processors. It is also not necessary that each processor group have four members, or that all processor groups have the same number of processors.

5 CPUs 100-107 each have individual primary caches 108. Typically, caches 108 will include at least separate data and instruction caches, each being, for example, 8K bytes of random access memory. In addition to the instruction and data cache components of caches 108, an additional cache memory of, for example, 1 megabyte of random access memory may be included as part of caches
10 108 in a typical system. CPUs 100-103 are connected to secondary cache 110 and make up a NUMA block and locale. CPUs 104-107 are connected to secondary cache 111. Caches 110 and 111 are connected via main system BUS 130 to shared memory 120. I/O BUS 140 connects disc array 125 and other typical data processing elements not shown. In the illustrated eight-processor embodiment,
15 secondary caches 110 and 111 are each 32 megabytes and shared memory 120 is 1 gigabyte of random access memory. Other sizes for each of these memory elements could readily have been employed.

In a system such as described, a "thread group" is a set of closely-related threads within a process that will tend to access and operate on the same data.
20 Handling these related threads as a single globally schedulable group promotes a closer relationship between the threads in the group and individual processors or groups of processors. This improves the ratio of cache hits and overall system performance.

The thread group is the basic unit of global scheduling across the system.
25 The thread group structure maintains the global scheduling policy and global priority which are used to schedule the execution of the thread group. The thread group structure also maintains the cumulative time slice and CPU accounting for all threads in its thread group, so time slicing and CPU accounting records for individual threads within a thread group are not necessary. Each individual thread
30 within the thread group maintains the thread priority and scheduling priority for itself.

Execution of a process often involves a plurality of thread groups, each with a plurality of threads. The use of thread groups in developing a process is well known from the prior art, for example, as described in detail in the previously-referenced U.S. Patent.

5 Looking again at Fig. 1, the various memory components within system 150 can be conceptualized as comprising three processing levels. Each level can be considered to contain one or more "instances" or nodes in the CPU/cache/shared memory hierarchy. Level 0 contains eight level 0 instances, or in the case of four locales or blocks, sixteen level 0 instances, each instance being
10 one of the CPUs 100-107 and its associated caches 108. Level 1 contains two level 1 instances, each comprising four level 0 instances, and a secondary cache. While eight processors are shown, the system can be arranged with more processors, for example, as a sixteen processor system. In the case of a sixteen processor system, level 1 contains four level 1 instances, each confirming four
15 level 0 instances and a secondary cache. Finally, there is a single level 2 instance containing the two level 1, or as applicable, the four level 1, instances and the shared system memory.

In the system, although the scheduling of threads is now well known, it has to now not been known how to keep all of the processors on a multi-processor
20 system performing productive work. As is well known from the prior art, each processor has a queue of jobs to perform. When the queue for a particular processor is empty, the processor goes idle. There may be times when one or more processors are idle, but the job queue of another processor is not empty. In such a case, it is necessary to decide which idle processor, if any, should receive
25 the work and how long it should wait before taking the job, i.e., poaching the process.

More specifically, a job will take the least amount of time if it is always run on the same processor. This is because of the multiple levels of memory caches used to reduce memory latency. In general, the lower the cache level, the less time
30 it takes to access the memory. When a job is run on a processor, its working set of memory is loaded into the applicable cache 108 of one of processors 100-107. On

005250" 03T 4960

a NUMA system, there is also another level of cache 111 that is shared between the four processors 100-103 or 104-107 of a functional block. In the best case scenario, the job is run on the processor which has the required code and data in its cache. The next best case scenario is to run the job on another processor on the same NUMA block. The tradeoff in this circumstance is between waiting for the processor on which the job last ran becoming free, or running the job on another idle processor. If time is allowed for the last processor, the idle processor's cycles will be wasted. If the job is run on a different processor, time and bus bandwidth may be wasted filling the cache of the idle processor and pulling all of the modified memory from the cache of the original processor. As noted previously, although the system and method are being described with reference to two blocks or locales 0 and 1, it will be appreciated that it can be implemented with, for example, sixteen CPUs, namely CPUs 0-3 in a locale 0, CPUs 4-7 in a locale 1, CPUs 8-11 in a locale 2, and CPUs 12-15 in a locale 3.

Turning now to the specific disclosure in the drawing, each processor 100-107 keeps track of the time from when it first went busy. It is intended that a process is poached when the processor on which it is scheduled is too busy to run the process. Accordingly, once a predetermined amount of time has elapsed during which a processor continues to be busy, other processors who have gone idle are then free to pouch the process from the busy processor. Each processor 100-107 has a timer for this purpose.

In accordance with the system and method, there is provided a strict ordering in the poaching order. Processors are paired up in a tree and start poaching from its nearest neighbor in the tree, i.e., electrically closest. Thus, in the case of Fig. 1, CPU0 on a first locale 0 (processors 100-103) will first try to poach from CPU1, then from CPU2, and CPU3, then from the locale 1, i.e., CPUs 104-107, and in the case of a sixteen-processor system, a CPU 15 on a locale 3 (not shown) will first try to poach from a CPU 14 (also not shown) and then from CPU 13 and CPU 12. This avoids the randomness that the current systems employ in poaching and tends to order poaching in pairs of electrically closest processors. In accordance with the system and method herein, as noted, each processor has an

independent clock that can be synchronized and eliminates wild poaching and provides the quickest response time.

The timing delay variables provided allow a system administrator to control and select a compromise between two extremes. If a delay is set high, it gives the user performance a boost on an idle system since all of the processes are effectively hard affined to their current locales. On the other hand, this means that there is a large system imbalance, and the scheduler will not straighten it out. If the delay is set small, the system load will be quickly balanced among the processors, but the overall throughput may actually degrade due to the cache fill overhead. The appropriate delay is a factor of the working set size of a typical process, the size of the caches, and the memory latency. Thus, the system administrator can then select the appropriate delay variable which maximizes through-put on the system.

Making reference to Fig. 2, the specifics of how poaching occurs is shown therein. More specifically, the process starts at a step 151 when a thread releases a processor and it goes idle. At step 153 the processor checks its level 0 run queue, level 1 run queue and level 2 run queue. At step 155 a determination is made about whether there is a thread group available on any one of these queues. If the answer is yes, then the process proceeds to immediately execute the highest priority thread through connection 157 to step 203 shown in Fig. 2C and discussed hereafter.

If the answer is no, at step 159, the processor increments level 0 idle count, level 1 idle count, marks level 0 and level 1 idle, and reads the current time. The idle indications are examined later by other idle processors looking for work. At step 161, every other run queue is sequentially checked and at step 163 a determination is made if this locale is marked to be skipped. The locale is marked to be skipped if there are any idle processors on it. In that case, it is always better to allow the idle processor from that same locale to pick up any work there.

If the locale is marked to be skipped, then at step 165 the next run queue is selected. If the answer is no, then at step 169 a determination is made about

whether the run queue next due time is less than the current time. If the answer is no, the process returns to step 165 as previously described, and if the answer is yes, at 171 connects to 171 in Fig. 2B. At step 173 the run queue next due time is recomputed and set equal to the time the run queue went busy plus the relative
5 time out. Note that this value is stored local to the processor doing the calculation so that it can be quickly checked later without incurring any bus traffic.

At step 175 a determination is made about whether the run queue next due time is less than the current time. If the answer is no, a connection is made at 177 to the process as further illustrated in Fig. 2D. If the answer is yes, at step 179 a
10 determination is made about whether there is work to do on a run queue. If the answer is no, the process is passed at 183 through the steps illustrated in Fig. 2E. If the answer is yes, at step 185 the highest priority thread in the run queue is selected, and the process continues at 187 as connected to Fig. 2C, where at step 189 the process decrements level 0 "is idle" count, and sets level 0 went busy time
15 to the current time.

At step 191, decrement level 1 is idle count, and at step 193 a determination is made about whether the level 1 idle count is equal to zero. If not equal to zero, the process is passed to step 197 discussed hereafter. If equal to zero, the process is passed to step 195 in which level 1 went busy time is set to the
20 current time, and at step 197 decrement level 2 "is idle" count.

At step 199 a determination is made if level 2 idle count is equal to zero. If not equal to zero, the process is passed to step 203 discussed hereafter. If equal to zero, at step 201, level 2 went busy time is set to the current time, and at step 203, the highest priority thread on a selected run queue is executed.

25 As may be appreciated from the previous discussion with reference to Fig. 2A, in addition to the above-discussed process, if at step 155 a determination is made that a thread group is available, the process passes directly to step 203 shown in Fig. 2C for execution of the highest priority thread, and after the highest priority thread is executed, then the process returns to the start at step 151.

Referring now to Fig. 2D, if at step 175 shown in Fig. 2B it is determined that the run queue next due time is not less than the current time, the process passes to step 205 where it checks if the run queue has any idle processors on its tree. If at step 207 it is determined that there are idle processors, then the locale is
5 marked to be skipped. If there are no idle processors, the process passes to 167 and return to step 165 previously discussed and shown in Fig. 2A.

Referring again to Fig. 2B, if at step 181 it is determined that there is no work to be done, the process passes to step 211 in Fig. 2E where the run queue next due time is set to equal the run queue time when the processor went busy plus
10 the relative time out, and then returns to step 205 in Fig. 2D.

Fig. 3 shows the process as the system is initialized and how the time periods are set.

At step 253 all processors are sequentially processed. Thereafter at step 255 the relatives tables are allocated. At step 257 each relative run queue is
15 processed in a manner such that at step 259 an entry is created in the relatives table containing a pointer to run queue, next due time, relative time out and locale index. The relative time out for the relative run queue is computed based on the locale number XOR relative locale number, the processor number XOR relative processor number, level 2 scheduling delay, level 1 scheduling delay, and level 0
20 scheduling delay. At step 261, the relatives table is sorted such that the shortest relative time out is at the top.

It will be appreciated from the above description that the invention may be implemented in other specific forms without departing from the spirit or essential characteristics thereof. The scope of the invention is indicated by the appended
25 Claims rather than by the foregoing description and all changes within the meaning and range of equivalency of the claims are intended to be embraced therein.

Claims

1. A method of allocating resources in a plurality of processors system, each processor of said plurality having a cache associated therewith, comprising:
- when at least one processor of said plurality of processors system is
- 5 idle, determining when at least one other processor of said plurality of processors is not idle;
- timing a predetermined period of time during which said at least one other processor which is not idle remains not idle; and
- if said at least one other processor which is not idle remains not idle
- 10 when the predetermined period of time has elapsed, poaching a process on a queue of said at least one non-idle processor to be run by said at least one processor which is idle.
2. The method of Claim 1 wherein if more than one processor is idle, poaching the process with the idle processor which is electrically closest to the at
- 15 least one non-idle processor.
3. The method of Claim 1 wherein the time period during which a non-idle processor is allowed to remain non-idle is greater the farther away a non-idle processor is electrically located relative to an idle processor.
4. The method of Claim 1 wherein an idle processor will first try to
- 20 poach from a non-idle processor electrically closest to it, and if the processor electrically closest to the idle processor is idle, it will then try to poach from the next processor which is electrically closest to it until it encounters a non-idle processor on which poaching can occur.
5. The method of Claim 1 wherein it is conducted on a ccNUMA
- 25 system.
6. The method of Claim 1 wherein each processor starts a timer associated therewith when it goes non-idle, and allowing an idle processor to poach a process therefrom only when a predetermined amount of time has elapsed on the timer.

7. The method of Claim 6 wherein said predetermined amount of time is established in relation to electrical proximity between an idle processor and a non-idle processor.

8. The method of Claim 7 wherein the greater the electrical distance
5 between idle and non-idle processors, the greater the predetermined amount of time which is allowed to elapse for said non-idle processor.

9. The method of Claim 1 wherein said processor system is a four block system comprised of sixteen processors arranged as CPU₀, CPU₁, CPU₂, CPU₃ on a locale 0, CPU₄, CPU₅, CPU₆, CPU₇ on a locale 1, CPU₈, CPU₉, CPU₁₀,
10 CPU₁₁ on a locale 2, and CPU₁₂, CPU₁₃, CPU₁₄, CPU₁₅ on a locale 3, and further comprising having each idle processor attempt to poach first from its nearest non-idle processor in its locale, and if there are no non-idle processors in said locale, from a processor in the closest locale having non-idle processors.

10. The method of Claim 9 wherein in the event more than one idle
15 processors attempt to poach a process from a non-idle processor, allowing the idle processor in closest electrical proximity to the non-idle processor poach the process.

11. A data processing system for allocating resources for simultaneously executing a plurality of processing tasks, comprising:

20 a plurality of processors, each having a cache associated therewith;

a timer associated with each processor for timing from the beginning of receiving a process from the processor's queue, the duration of time the process is run, during which time the processor is running the process is non-idle;

25 means for determining the duration of time any one processor is not idle;

means for poaching a process from a non-idle processor by an idle processor when said duration of time during which the non-idle processor is running a process exceeds a predetermined amount.

006250"0274960

12. The system of Claim 11 further comprising means for allowing an idle processor electrically closest connected to a non-idle processor to poach a process in the event there are more than one non-idle processors.
13. The system of Claim 1 further configured to allow the time period
5 during which a processor is allowed to remain non-idle is determined in accordance with proximity in electrical connection between the non-idle processor and a non-idle processor, wherein the greater the connection distance, the greater the predetermined amount of time.
14. The system of Claim 11 configured for having an idle processor first
10 try to poach from a non-idle processor electrically closest to it, and if the processor electrically closest to it is idle, then to try to poach from the next electrically closest processor until it encounters a non-idle processor on which poaching can occur.
15. The system of Claim 11 wherein said system is a ccNUMA system.
16. The system of Claim 11 wherein each processor is configured for
15 starting the timer associated therein when it goes non-idle.
17. The system of Claim 16 wherein each processor is configured for setting the predetermined amount of time in relation to electrical connection proximity between an idle processor and a non-idle processor.
18. The system of Claim 17 wherein each processor is configured for
20 setting the predetermined amount of time such that the greater the electrical connection distance between an idle and a non-idle processor, the greater the amount of time which is allowed to elapse before an idle processor is allowed to poach a process from a non-idle processor.
19. The system of Claim 1 wherein the plurality of processors are
25 arranged in four blocks comprised of sixteen processors arranged as CPU₀, CPU₁, CPU₂ and CPU₃ on a locale 0, CPU₄, CPU₅, CPU₆ and CPU₇ on a locale 1, CPU₈, CPU₉, CPU₁₀ and CPU₁₁ on a locale 2, and CPU₁₂, CPU₁₃, CPU₁₄ and CPU₁₅ on a locale 3, and comprising said processors configured such that each idle processor
30 attempts to poach first from its nearest non-idle processor in its locale, and if there

are no non-idle processors in said locale, from a processor in the closest locale having non-idle processors.

20. The system of Claim 19 wherein said processors are configured for allowing an idle processor in closest electrical connection to a non-idle processor
- 5 to poach a process therefrom in the event more than one idle processor attempts to poach a process from said non-idle processor.

006250" 0274960

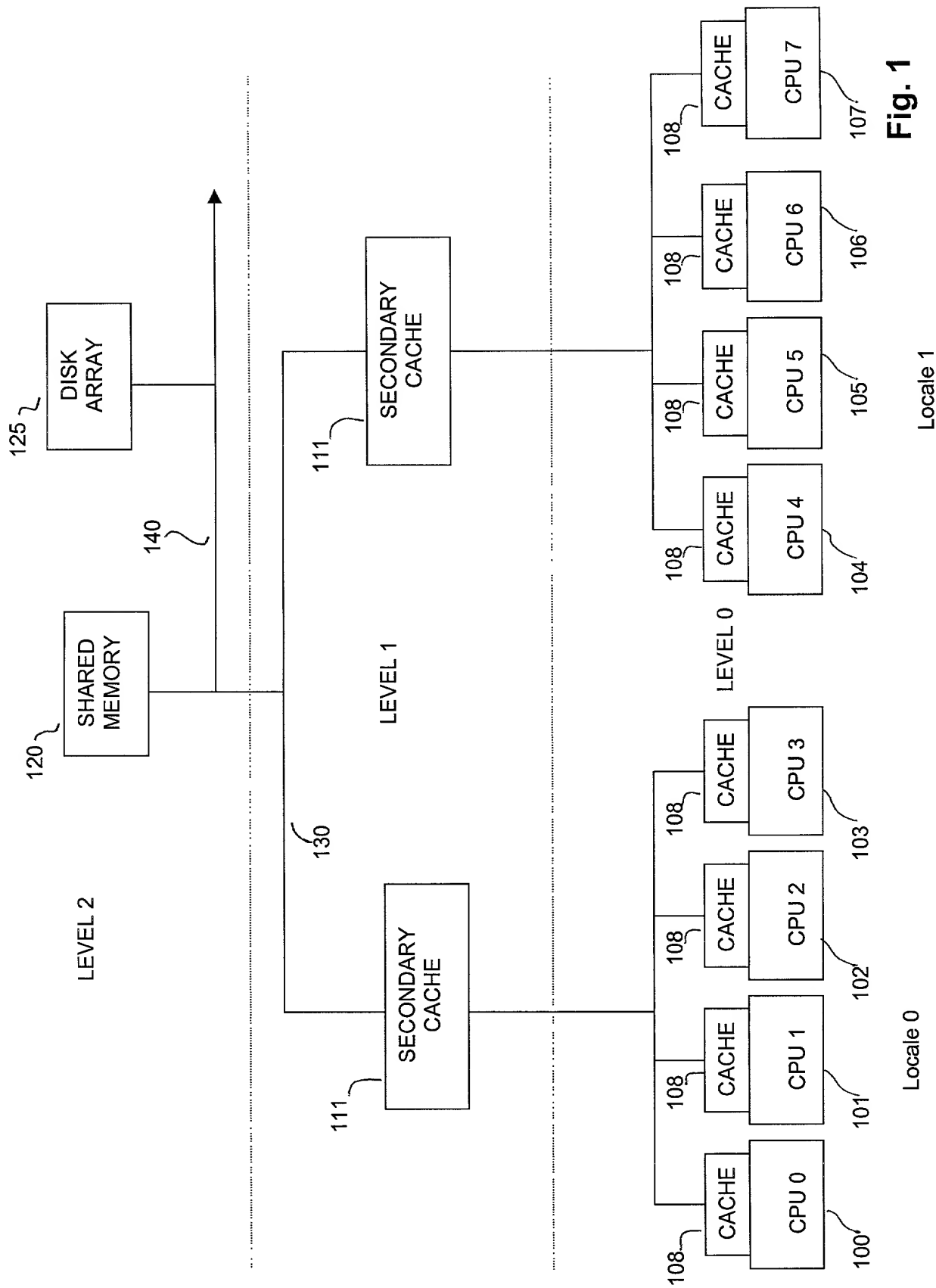
Abstract

A method and system allocate resources in a plurality of processors system. When a processor is idle, the system determines when another processor is not idle. The time the non-idle processor remains non-idle is timed, and once a
5 predetermined amount of time elapses, if the non-idle processor is still not idle, the idle processor poaches a job from the non-idle processor.

RALLIB01:570543.2

10

005260" 027.2360



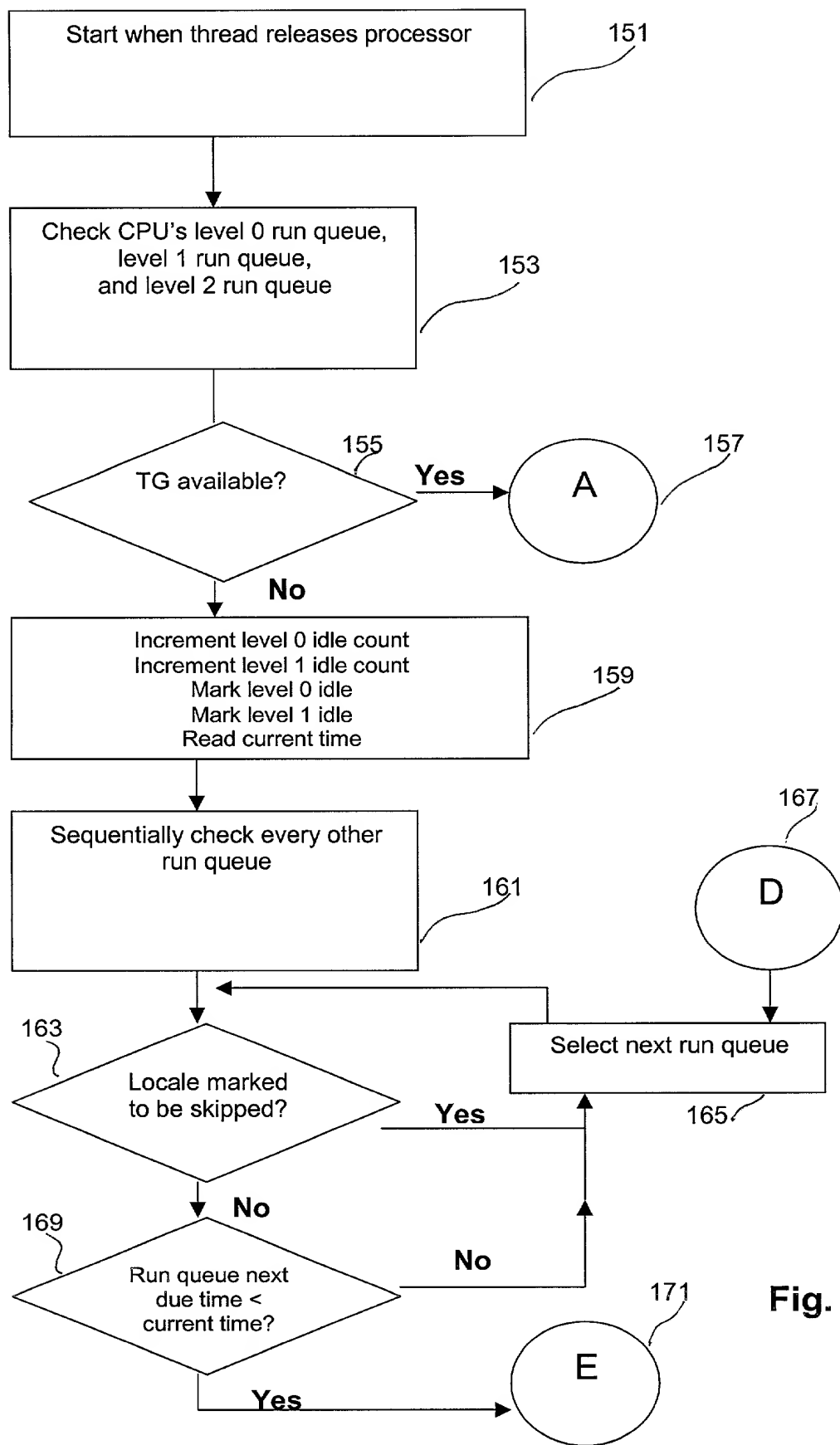


Fig. 2A

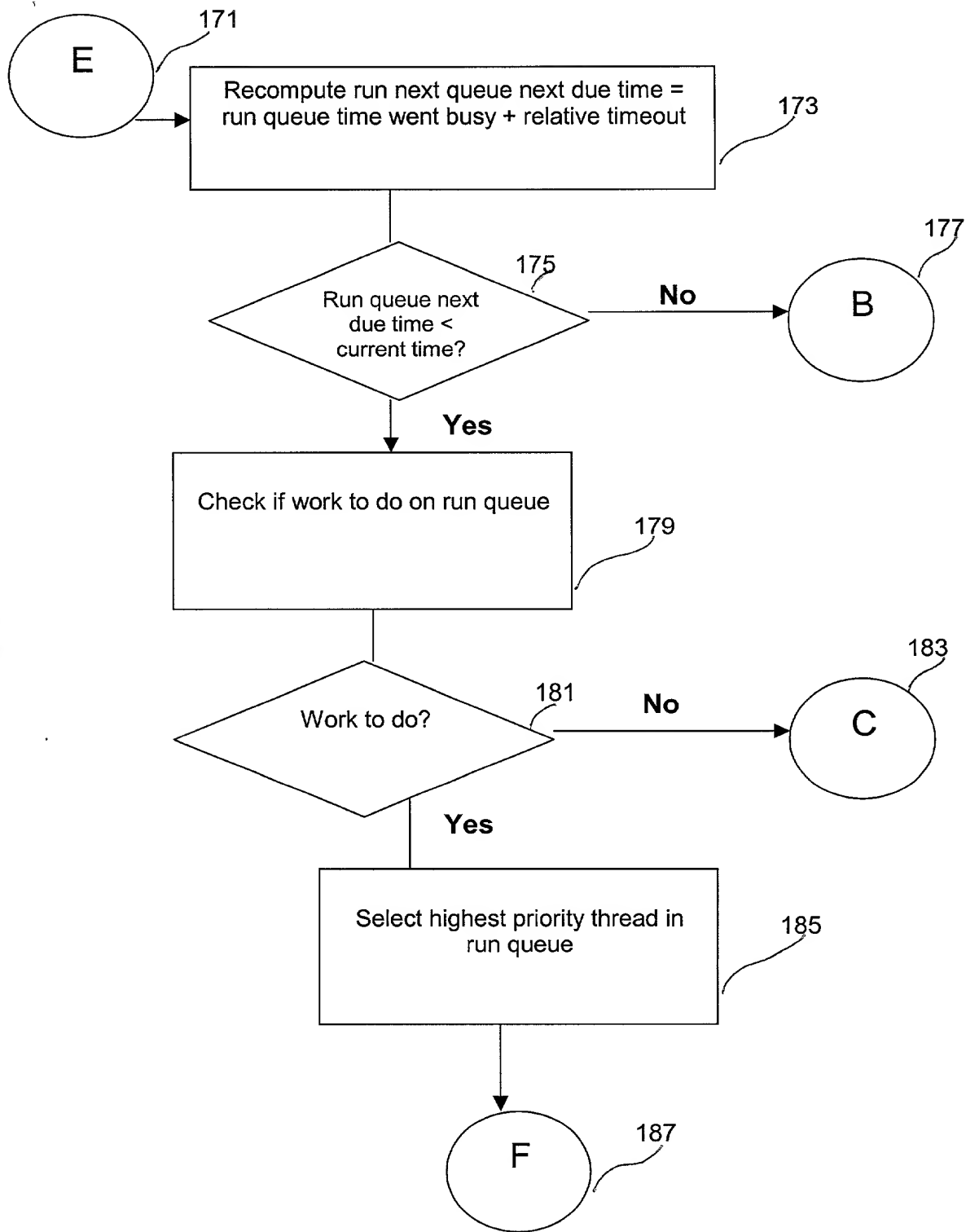


Fig. 2B

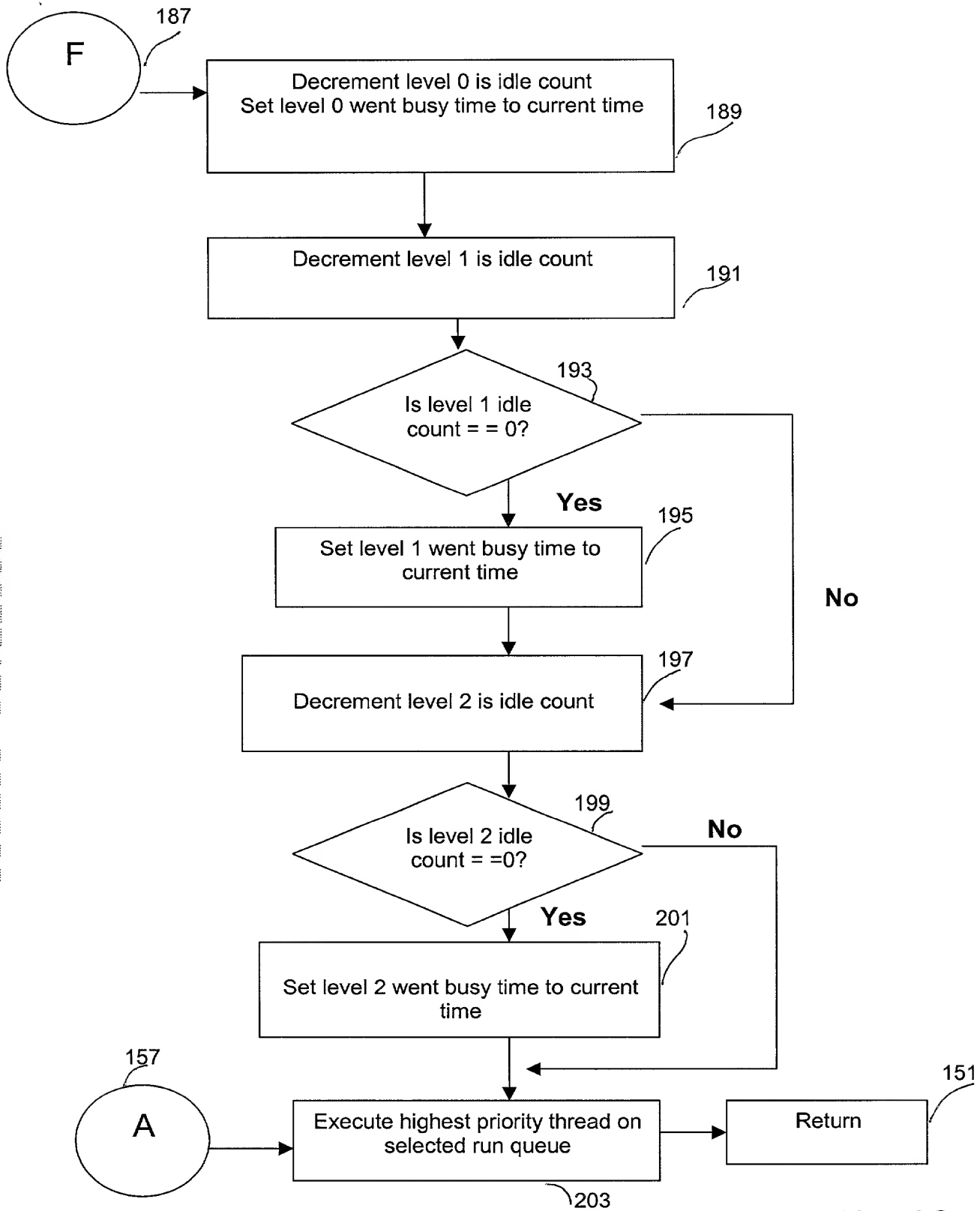


Fig. 2C

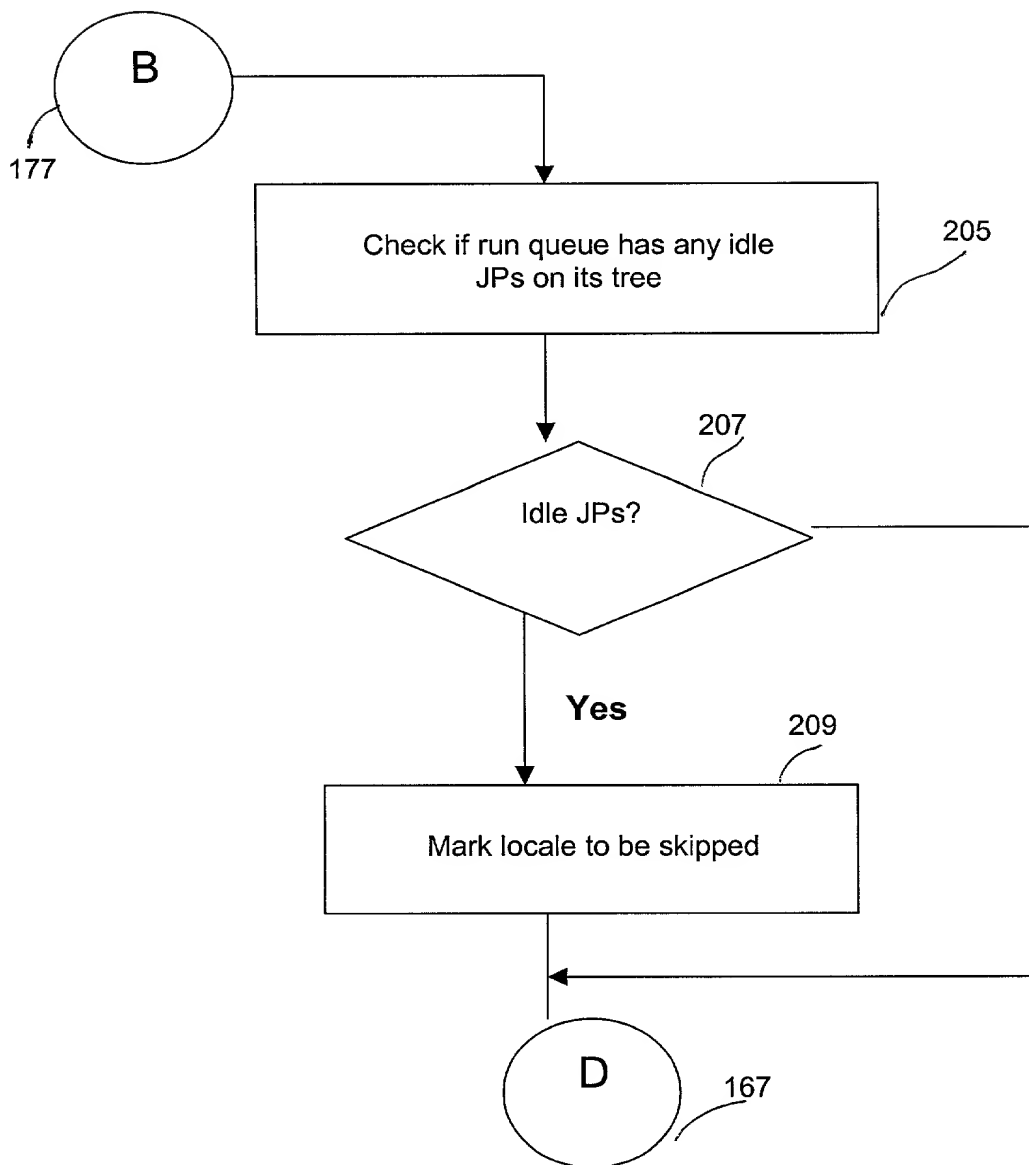


Fig. 2D

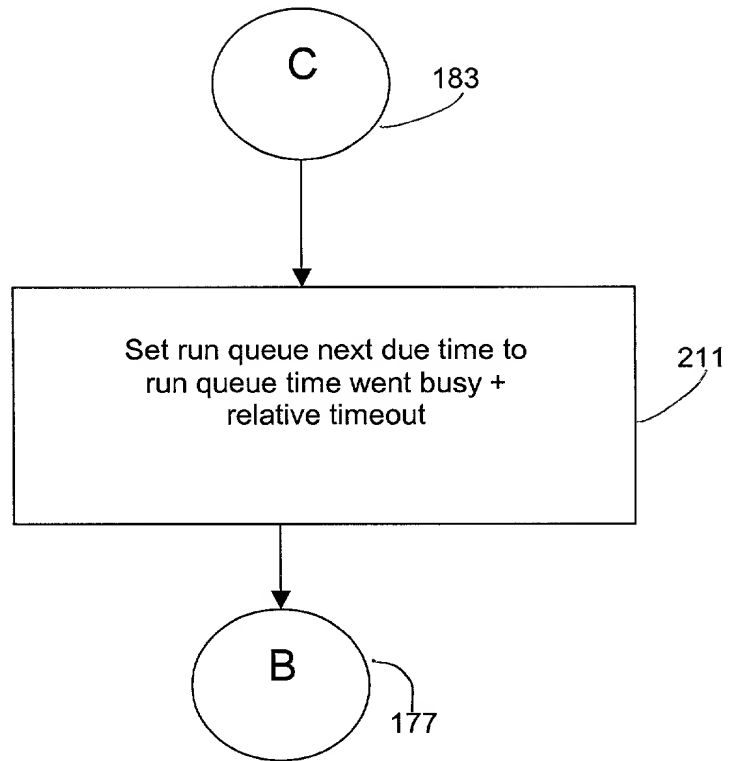


Fig. 2E

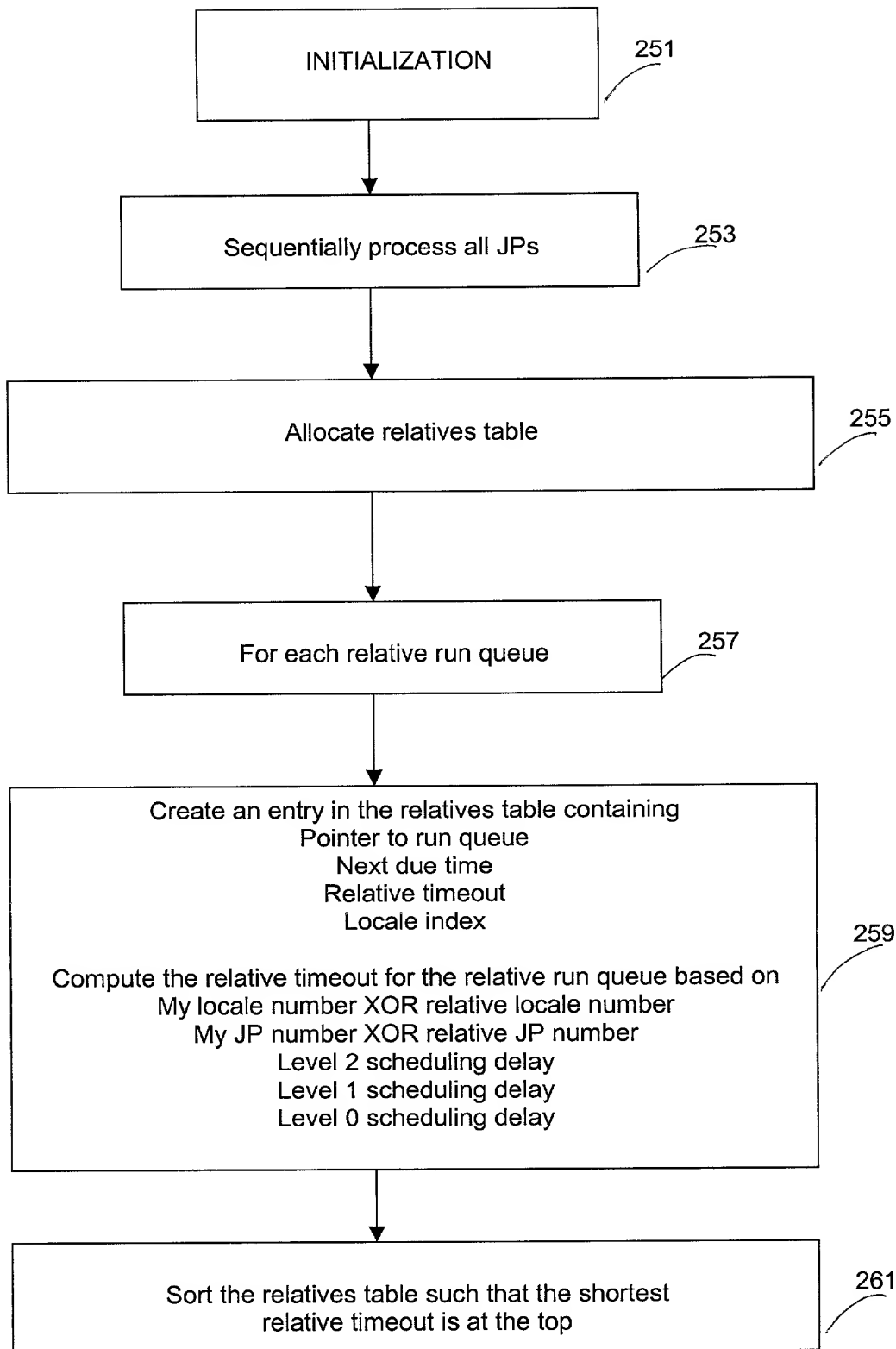


Fig.3

DECLARATION AND POWER OF ATTORNEY FOR PATENT APPLICATION
English Language Declaration

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name.

I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled **APPARATUS AND METHOD FOR PROCESS DISPATCHING BETWEEN INDIVIDUAL PROCESSORS OF A MULTI-PROCESSOR SYSTEM**; the specification of which: (check one)

X is attached hereto.

_____ was filed on _____ as

Application Serial No. _____

_____ and was amended on _____ (if applicable)

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose information which is material to patentability as defined in Title 37, Code of Federal Regulations, §1.56.

I hereby claim foreign priority benefits under Title 35, United States Code, §119 of any foreign application(s) for patent of inventor's certificate listed below and have also identified below any foreign application for patent or inventor's certificate having a filing date before that of the application on which priority is claimed:

Prior Foreign Application(s)

Priority Claimed

NONE	
(Number)	(Country)

(Day/Month/Year Filed)

<u>Yes</u>	<u>No</u>
------------	-----------

(Number)	(Country)
----------	-----------

(Day/Month/Year Filed)

<u>Yes</u>	<u>No</u>
------------	-----------

I hereby claim the benefit under Title 35, United States Code, §120 of any United States application(s) listed below and insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States application in the manner provided by the first paragraph of Title 35, United States Code §112, I acknowledge the duty to disclose material to patentability as defined in Title 37, Code of Federal Regulations, §1.56 which became available between the filing date of the prior application and the national or PCT international filing date of this application:

(Application Serial No.)

(Filing Date)

(Status)
(patented, pending, abandoned)

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true, and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

English Language Declaration

POWER OF ATTORNEY: As a named inventor, I hereby appoint the following attorney(s) and/or agent(s) to prosecute this application and transact all business in the Patent and Trademark Office connected therewith.

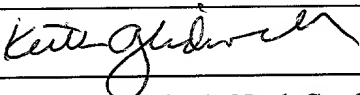

A. José Cortina, Registration No. 29,733; George T. Marcou, Registration No. 33,014; Benjamin Driscoll, Registration No. 41,571, Charles W. Calkins, Registration No. 31,814, John M. Harrington, Registration No. 25,592, Dawn-Marie Bey, Registration No. 44,442, John Ball, Registration No. 44,433; James J. Bindseil, Registration No. 42, 326; J. Steven Gardner, Registration No. 41,772; and Jason Link, Registration No. 44,874.

Send Correspondence to:

A. José Cortina
Kilpatrick Stockton LLP
3737 Glenwood Avenue, Suite 400
Raleigh, North Carolina 27612

Direct telephone calls to:

A. José Cortina
(919) 420-1820

Full name of sole inventor: Keith Glidewell	
Inventor's Signature 	Date: September <u>27</u> , 2000
Residence: 3328 Boulder Court, Raleigh, North Carolina 27607	
Citizenship: United States of America	
Post Office Address: 3328 Boulder Court, Raleigh, North Carolina 27607	
	

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re application of: Keith Glidewell
Application No.: Unassigned
Filed: Herewith
For: APPARATUS AND METHOD FOR PROCESS
DISPATCHING BETWEEN INDIVIDUAL PROCESSORS
OF A MULTI-PROCESSOR SYSTEM

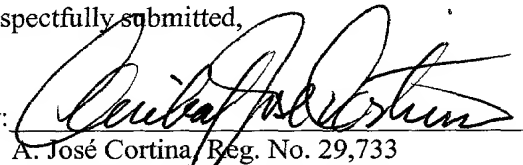
ASSOCIATE POWER OF ATTORNEY

The Honorable Assistant Secretary and
Commissioner of Patents and Trademarks
Box Patent Application
Washington, D.C. 20231

I, Anibal José Cortina, Attorney of Record in the above-identified application, hereby grant to Robert Dulaney, Reg. No. 28071, an Associate Power of Attorney, with full power of substitution and revocation to prosecute and transact all business in the Patent and Trademark Office connected therewith.

Respectfully submitted,

By:




A. José Cortina, Reg. No. 29,733
One of the Attorneys for Applicant
Kilpatrick Stockton LLP
3737 Glenwood Avenue, Suite 400
Raleigh, North Carolina 27612
Telephone: (919) 420-1820

Date: September 29, 2000

Certificate of Mailing

I hereby certify that this document is being deposited with the United States Postal Service as Express Mail in an envelope addressed to the Assistant Commissioner for Patents, Box Patent Application, Washington, D.C. 20231 on September 29, 2000, Express Mail Label No. EJ841202327US.


Karyn Grinn, Assistant
RALL1501:571908.1

006677120-092900